1

2

# Using machine learning models to predict oxygen saturation following ventilator support adjustment in critically ill children: a single center pilot study.

6

7

8  Sam Ghazal[1], Michael Sauthier[2], David Brossier[2], Wassim Bouachir[3], Philippe Jouvet[2*] and
9  Rita Noumeir[1]

10

11

12  [1] Department of health information analysis, École de Technologie Supérieure (ÉTS),

13  Montreal, Quebec, Canada

14  [2] Department of Pediatrics, Sainte-Justine Hospital, Montreal, Quebec, Canada

15  [3] LICEF research center, TÉLUQ University, Montreal, Quebec, Canada

16

17

18  * Corresponding author:

19  philippe.jouvet@umontreal.ca (PJ)

# Abstract

Clinicians' experts in mechanical ventilation are not continuously at each patient's bedside in an intensive care unit to adjust mechanical ventilation settings and to analyze the impact of ventilator settings adjustments on gas exchange. The development of clinical decision support systems analyzing patients' data in real time offers an opportunity to fill this gap. The objective of this study was to determine whether a machine learning predictive model could be trained on a set of clinical data and used to predict hemoglobin oxygen saturation 5 min after a ventilator setting change. Data of mechanically ventilated children admitted between May 2015 and April 2017 were included and extracted from a high-resolution research database. More than $7.10^5$ rows of data were obtained from 610 patients, discretized into 3 class labels. Due to data imbalance, four different data balancing process were applied and two machine learning models (artificial neural network and Bootstrap aggregation of complex decision trees) were trained and tested on these four different balanced datasets. The best model predicted $SpO_2$ with accuracies of 76%, 62% and 96% for the $SpO_2$ class "< 84%", "85 to 91%" and "> 92%", respectively. This pilot study using machine learning predictive model resulted in an algorithm with good accuracy. To obtain a robust algorithm, more data are needed, suggesting the need of multicenter pediatric intensive care high resolution databases.

2

# Introduction

In case of respiratory failure, mechanical ventilation supports the oxygen ($O_2$) diffusion into the lungs and the carbon dioxide ($CO_2$) body removal. As an expert in mechanical ventilation cannot reasonably be expected to be continuously present at the patient's bedside, specific medical devices aimed to help in ventilator settings adjustments may help to improve the quality of care. Such devices are developed using either algorithms based on respiratory physiology/medical knowledge that adapt ventilator settings in real time based on patients' characteristics but are not accurate enough to be used widely in clinical practice, especially in children [1, 2]; or physiologic models that simulate cardiorespiratory responses to mechanical ventilation settings modifications but none was validated for this indication [3]. The above-mentioned models all share the limitation of not being suited to learn from ever-growing sets of clinical research data, and potentially improve their performances. To overcome this drawback, another avenue is the development of algorithms using artificial Intelligence to provide caregivers with support in their decision-making tasks. In this study, we assessed machine learning methods to predict transcutaneous hemoglobin saturation oxygen ($SpO_2$) of mechanically ventilated children after a ventilator setting change using a high resolution research database.

# Materials and Methods

This study was conducted at Sainte-Justine Hospital and included the data collected prospectively between May 2015 and April 2017 of all the children, age under 18 years old, admitted to the Pediatric Intensive Care Unit (PICU) who were mechanically ventilated with an endotracheal tube. Patients' data were excluded if the patient was hemodynamically unstable defined as 2 or more vasoactive drugs delivered at the same time (ie., epinephrine,

3

63    norepinephrine, dopamine or vasopressin) or with an uncorrected cyanotic heart disease

64    defined by no $SpO_2 > 97\%$ during all PICU stay. All the respiratory data from included

65    patients were extracted from the PICU research database [4], after study approval by the

66    ethical review board of Sainte-Justine hospital (number 2017 1480).

67

## 68    Data extraction

69    To determine the data that will be extracted for each child, an item generation was

70    conducted by three physicians (PJ, MS, DB). The resulting items are presented in Fig 1 within

71    their sources, means of extraction and a schematic of the main components of the study.

72    The predictive $SpO_2$ value was the $SpO_2$ 5 minutes after a change of a ventilator setting. The

73    delay of 5 min corresponded to the shortest period of time to reach a steady state after

74    modification of a ventilator setting [5].

75    **Fig 1. Schematic description of the analysis process and items involved.** EMR: electronic
76    Medical Record, $FiO_2$: inspired fraction of Oxygen, Vt: tidal volume, PEEP:
77    Positive end expiratory pressure, PS above PEEP: pressure support level Above
78    PEEP, PC above PEEP: pressure control level above PEEP, MVe: expiratory
79    minute volume, I:E Ratio: inspiratory time over expiratory time, Measured RR:
80    respiratory rate measured by the ventilator, PIP: positive inspiratory pressure ie
81    maximal pressure measured during inspiration. $_{5min}SpO_2$: $SpO_2$ observed 5 min
82    after PEEP, $FiO_2$, tidal volume, PS above PEEP, PC above PEEP change, ML:
83    machine learning, ANN: artificial neural network, BACDT: Bootstrap aggregation
84    complex decision trees.
85

## Data formatting

The data extracted from the research database needed: (1) to remove erroneous data due to disconnection of the patient from the ventilator or the monitor, or due to transient interventions such as suctioning; (2) to remove the rows at which no ventilator setting variables was modified; (3) to adapt data format for classifier training. The methodology to format the data is described in S1 file.

## Data categorization

$SpO_2$ levels at 5min were classified into three categories (Table 1). The thresholds were selected according to clinical value: a $SpO_2 < 92\%$ is a target to increase oxygenation in mechanically ventilated children [6]. The critical level of 85% $SpO_2$ is used as an alarm of severe hypoxemia in intensive care [7].

**Table 1: Definition of SpO$_2$ class labels specifications**

| $SpO_2$ classification | $SpO_2$ range (%) | Rows number (n) |
|---|---|---|
| 1 | < 84 | 17,112 |
| 2 | 85 to 91 | 29,869 |
| 3 | 92 to 100 | 729,746 |

## Data balancing

The data analysis showed a severe imbalance with most $SpO_2$ at 5min above 92%. This is logical as caregivers want to maintain $SpO_2$ in normal range during child PICU stay. In such condition, the classifier learns the majority class label (class 3) (Table 1) but doesn't learn the minority class labels (class 1 and 2) [8]. The data balancing process aims to allow the

5

106    classifier to learn from all class equally. The data balancing process used in this study

107    included a combination of down-sampling and up-sampling techniques: to balance the three

108    classes of the data involved, a down-sampling of the $SpO_2$ class 3 using TOMEK algorithm [9]

109    and an over-sampling of $SpO_2$ class 1 and 2 using Synthetic Minority Oversampling

110    Technique (SMOTE) [10] were performed.

111    The creation of synthetic data points by SMOTE can be formulated as follows:

112    $$x_{syn} = x_i + (x_{knn} - x_i) \times \delta \tag{2}$$

113    In equation (2), $x_{syn}$ represents the synthetic data point. The variables $x_i$ and $x_{knn}$ are

114    respectively the original instance, and the nearest neighbor data point which is randomly

115    picked among the $k$ nearest neighbors.. The random number $\delta$ is generated in [0,1] to

116    determine the position of the created synthetic data point along a straight line joining the

117    original data point $x_i$ and its chosen nearest neighbor $x_{knn}$.

118    To study which data balancing method provided the more accurate algorithm, four datasets

119    were produced via four different balancing procedures, involving different combinations of

120    data balancing techniques (Fig 2).

121    **Fig 2. Descriptions of the four balancing procedures.**

122

# Predicted $SpO_2$ Classification

124    To identify the best machine learning classification method, we tested two classification

125    models: artificial neural network and bagged complex decision trees, on the four balanced

126    datasets.

## Artificial Neural Network (ANN)

128    Once the data has been pre-processed, a machine learning predictive model was trained on a

129    sub-set of labeled training data. The model is then used to predict the target variable values on

130    a testing subset where the class labels are hidden.  We used Artificial Neural Networks

131    (ANN) to make predictions of the $SpO_2$ variable, based on the values of other variables of

132    interest. Through the function approximation that the ANN performs, it is possible to make

133    predictions of $SpO_2$ variable, based on the input data.

134

135    The ANN is learned from training data, using the backpropagation algorithm [11] and is

136    tested on a test set made of the remaining rows of data to validate the generalization of the

137    model. The learning algorithm runs through all the rows of data in the training data set and

138    compares the predicted outputs with the target outputs found in the training data set. The

139    weights are adjusted via supervised learning, in a manner to minimize the error of predicted

140    $SpO_2$ vs target $SpO_2$. The process is repeated until the error is minimized.

141

142    The ANN classifier was implemented through cycles of forward propagation followed by

143    backward propagation through the network's layers.  The backpropagation algorithm is used

144    for performance optimization.

145    For a given number of classes K > 2, the cross-entropy error can be formulated as shown in

146    eq. 3, where $(W_i)_i$ is the matrix of weights between the neuron layers, $r_i$ is the target value. $y_i$ is

147    the value generated by the ANN, ie., its output.

148
$$E^t((W_i)_i \mid x^t, r^t) \;=\; - \sum_i r_i^t \log y_i^t \qquad (3)$$

149     The outputs of the ANN are:

150
$$y_i^t = \frac{\exp w_i^t x^i}{\sum_k \exp w_k^t x^t} \qquad (4)$$

151    Using stochastic gradient-descent (SGD) for error minimization, the update rule for the ANN

152    weights is:

153
$$\Delta w_{ij}^t = n(r_i^t - y_i^t)x_j^t \qquad (5)$$

7

154     In equation 5, $\eta$ is the learning rate, which, when SGD is used, decreases as the error is

155     minimized. During ANN training, each observation, comprised of an input vector and a target

156     output, is denoted ($x^t$, $r^t$), with $r^t \in$ ("1", "2", "3"). The reason why the cross-entropy (eq. 3) is

157     used instead of the Least Square Error (LSE) is to avoid long periods of training, due to the

158     ANN going through stages of slow error reduction.

159

## Bootstrap aggregation of complex decision trees

161     Bootstrap aggregating (acronym: bagging) was proposed by L Breiman in 1994 to improve

162     classification by combining classifications of randomly generated training sets [12]. Bagging

163     allows for the creation of an aggregated predictor via the use of multiple training sub-sets

164     taken from the same training set. Let ($T^i$) denote the replicate training sub-sets bootstrapped

165     from the training set $T$. These replicate sub-sets each contain $N$ observations, drawn at

166     random and with replacement from $T$. For each of these sub-sets of $N$ observations, a

167     prediction model, or classifier, is created. The computational model we used for bagging was

168     complex decision trees. This means that, for each bootstrapped sub-set of training data, a

169     complex decision tree is trained and thus a classifier is created. If $i = 1, ..., n$, then $n$

170     classifiers are created through the bagging process.

171

172     A decision tree is a flowchart computational model which can be used for both regression, as

173     well as classification problems. Paths from the root of the tree to its various leaf nodes go

174     through decision nodes in which decision rules are applied in a recursive manner, based on

175     values of input variables. Each path represents an observation $(X, y) = (x_1, x_2, x_3, ..., x_n, y)$,

176     where the label assigned to the target $y$ is given in the leaf node, at the end of the path, ie.,

177     classification [13].

178

179 In the aim of maximizing the model's generalization capability during the training process,

180 the Bagged Complex Trees' performance is tested via $k$-fold cross-validation. A value $k = 10$,

181 which is common practice, was used in this study. The training using $k$-fold cross-validation

182 is carried out as described in Fig 3.

183 **Fig 3. *k*-fold cross-validation**

184

185 The mathworks Matlab R2016b Machine Learning toolbox was used for the creation of the

186 ensemble of Bagged complex trees model.

187
## Assessment of the performances of the classifiers

189 We evaluated the performances of the classifiers based on the metrics including testing

190 confusion matrix, average accuracy, precision, recall and F score [14] with a $_{5min}SpO_2$

191 prediction expected above 0.9 for each class.

192

193 • Precision

194
$$Precision = \frac{\# \; True \; positives \; class \; i}{Total \; \# \; classifications \; for \; class \; i} \quad (6)$$

195 The *Precision* (eq. 6) is the ratio of all correct classifications for class *i* to all instances labeled

196 as class label *i* by the model. In a non-normalized confusion matrix, this would mean

197 dividing the number of instances classified in class label *i* by the total of instances in column

198 *i.*

199

200 ▪ Recall

201
$$Recall = \frac{\# \; True \; positives \; class \; i}{Total \; \# \; observations \; class \; i} \quad (7)$$

202 Recall is the ratio of the number of instances classified in class label $i$ to the number of true

203 class $i$ labels. In a non-normalized matrix, this would require dividing the number of

204 instances classified in class label $i$ by the total of row $i$

205

206 • F-score

207
$$F - score = \cfrac{2}{\cfrac{1}{recall} + \cfrac{1}{precision}} \quad (8)$$

208 The F-score provides a single measure of classification performance of the model used.

209

## 210 Results and discussion

211 We developed and assessed the performances of two machine learning classifiers on four

212 different balanced datasets to predict SpO$_2$ at 5 min after a ventilator setting change (*ie* FiO$_2$,

213 PEEP, Vt/Pressure), in 610 mechanically ventilated children. In Fig 4 and Table 2, we report

214 the performances of these two classifiers. Using the classification performance metrics, the

215 bagged trees classifier trained on dataset #3 (see Fig 2) has yielded the best classification

216 performance on the test sets (Table 2). The confusion matrix of the whole bagged trees

217 shows that SpO$_2$ at 5 min could correctly predict in **76%** of class "1" data, **62%** of class "2",

218 and **96%** of class "3" (Fig 4).  This huge variation in classification performances of the three

219 class labels can be explained by the large variation in the numbers of observations available

220 for each of the class labels in the initial dataset that has limited the machine learning (Table

221 1).

222

223 **Fig 4. Artificial neural network (ANN) and bootstrap aggregation of complex decision trees**
224 **(BACDT) test confusion matrices.** The darker colors represent higher levels of accuracy. A:

225  balanced dataset 1, B: balanced dataset 2, C: balanced dataset 3, D: balanced dataset 4 (see
226  Fig 2).
227

228  **Table 2. Performance of artificial neural networks (ANN) and bootstrap aggregation of**
229  **complex decision trees (BACDT) classifiers for SpO₂ prediction at 5 min following a**
230  **ventilator setting change.** Avg/total: average accuracy of total classification values. In italics
231  is the performance of the best predictive model obtained among the eight tested.
232

| Balanced datasets | $_{5min}SpO_2$ class | ANN | | | BACDT | | |
|---|---|---|---|---|---|---|---|
| | | Precision | Recall | F-score | Precision | Recall | F-score |
| **Dataset 1** | 1 | 0.12 | 0.70 | 0.21 | 0.80 | 0.76 | 0.78 |
| | 2 | 0.16 | 0.43 | 0.23 | 0.61 | 0.56 | 0.59 |
| | 3 | 0.96 | 0.67 | 0.79 | 0.97 | 0.98 | 0.97 |
| | Avg/total | 0.88 | 0.65 | 0.73 | 0.94 | 0.94 | 0.94 |
| **Dataset 2** | 1 | 0.09 | 0.72 | 0.16 | 0.77 | 0.72 | 0.74 |
| | 2 | 0.09 | 0.47 | 0.16 | 0.57 | 0.53 | 0.55 |
| | 3 | 0.98 | 0.70 | 0.81 | 0.98 | 0.99 | 0.98 |
| | Avg/total | 0.93 | 0.69 | 0.78 | 0.96 | 0.97 | 0.97 |
| **Dataset 3** | 1 | 0.16 | 0.68 | 0.25 | *0.80* | *0.76* | *0.78* |
| | 2 | 0.26 | 0.42 | 0.33 | *0.67* | *0.62* | *0.65* |
| | 3 | 0.92 | 0.60 | 0.72 | *0.95* | *0.96* | *0.96* |
| | Avg/total | 0.80 | 0.58 | 0.65 | *0.91* | *0.91* | *0.91* |
| **Dataset 4** | 1 | 0.09 | 0.69 | 0.16 | 0.80 | 0.74 | 0.77 |
| | 2 | 0.12 | 0.47 | 0.19 | 0.58 | 0.54 | 0.56 |
| | 3 | 0.97 | 0.68 | 0.80 | 0.98 | 0.98 | 0.98 |
| | Avg/total | 0.92 | 0.67 | 0.76 | 0.96 | 0.96 | 0.96 |

233

234  For the artificial neural network, the variation of the number of hidden layers and number of

235  neurons per hidden layer did not seem to have a significant effect on the model's

236  classification performance (Table 3). As for the Bagged complex trees, the variation of the

237  number of complex trees did not yield significant changes in classification performance

238  (Table 4).

239

240  **Table 3. Absence of impact on performance of the increase of neurons and hidden layers**
241  **for artificial neural network (ANN).** Example of the performance assessed by the F score on
242  the balanced dataset 3 (see fig 2)

243

244

| ANN | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Hidden layers (n) | | 1 | | | 2 | | | 3 | |
| Neurons/hidden layer (n) | | 10 | 50 | 100 | 10 | 50 | 100 | 10 | 50 | 100 |
| **F-score** | $_{5min}$SpO$_2$ class 1 | 25 | 25 | 25 | 25 | 25 | 25 | 22 | 22 | 19 |
| | $_{5min}$SpO$_2$ class 2 | 33 | 33 | 33 | 33 | 33 | 33 | 33 | 33 | 32 |
| | $_{5min}$SpO$_2$ class 3 | 72 | 72 | 72 | 72 | 72 | 72 | 69 | 69 | 69 |

245

246 **Table 4. Absence of impact on performance of the number of complex trees for bootstrap**
247 **aggregation of complex decision trees (BACDT).** Example of the performance assessed by
248 the F score on the balanced dataset 3 (see Fig 2)

249

| | | BACDT | |
|---|---|---|---|
| | | n = 30 | n=50 |
| **F-score** | $_{5min}$SpO$_2$ class 1 | 78 | 78 |
| | $_{5min}$SpO$_2$ class 2 | 65 | 65 |
| | $_{5min}$SpO$_2$ class 3 | 96 | 96 |

250

251 In agreement with previous studies regarding bagging being a better method for medical

252 data classification, tree Bagging fared better than the artificial neural network used in this

253 study [12]. It is noteworthy however that the gaps in performance results between the

254 training and testing confusion matrices are relatively higher in the case of bagged trees

255 model than in that of the artificial neural network (Fig 5). This seems to indicate that,

256 although the bagged trees model was capable of learning very well from the data, there's

257 still room for improvement in the generalization. The SMOTE algorithm is designed in such a

258 way that should theoretically not affect the generalization of the trained model. In cases of

259 extreme data imbalance, however, as is the case in this study, the over-sampling within the

260 data space of a given minority class label, used for increasing the cardinality of the class

261 label's set, is also likely to be extreme. This may render the data space of this class relatively

262 dense with respect to the rest of the data, made up of real data points of the studied patient

263 sub-population. This may potentially explain the classification model's relatively poor

264    generalization for $_{5min}SpO_2$ class "1" and "2" with respect to the generalization for $_{5min}SpO_2$

265    class "3". Also, since SMOTE generates synthetic data points by interpolating between

266    existing minority class instances, it can obviously increase the risk of over-fitting when

267    classifying minority class labels, since it may duplicate minority class instances.  The fact that

268    the training confusion matrix shows extremely high classification performances for the

269    minority $_{5min}SpO_2$ class "1" and "2", as opposed to those shown in the testing confusion

270    matrix, suggests that the over-sampling of the minority $_{5min}SpO_2$ class using SMOTE could

271    have caused some overfitting for these classes, but this would have to be further

272    investigated.

273

274    **Fig 5. Training and testing confusion matrices of artificial neural networks (ANN) and**
275    **bootstrap aggregation of complex decision trees (BACDT) classifiers for SpO$_2$ prediction at**
276    **5 min following a ventilator setting change.**
277

278    The strengths of this study include a large clinical database of mechanically ventilated

279    children used with more than $7.10^5$ rows. In a recent similar study in PICU, 200 patients were

280    included with $1.15.10^3$ rows [15]. However, the volume of data is clearly insufficient. To use

281    such machine learning predictive models, the pediatric intensive care community needs to

282    combine multicenter high resolution database. In addition, children data could be pooled to

283    neonatal and adult intensive care data, when possible, such as MIMIC III database [16]. The

284    other strength is the process used to transform the data into a usable format and to correct

285    a variety of artifacts present (S1 file). In health care, there is a significant interest in using

286    clinical databases including dynamic and patient-specific information into clinical decision

287    support algorithms. The ubiquitous monitoring of critical care units' patients has generated a

288    wealth of data which presents many opportunities in this domain. However, when

289     developing algorithms domains, such as transport or finance, data are specifically collected

290     for research purposes. This is not the case in healthcare where the primary objective of data

291     collection systems is to document clinical activity, resulting in several issues to address in

292     data collection, data validation and complex data analysis [17]. As detailed in S1 file, a

293     significant amount of effort is needed, when data have been successfully archived and

294     retrieved, to transform the data into a usable format for research.

295     This study has several limitations. The limited row number reduced the $SpO_2$ classification

296     for machine learning predictive model to three clinically relevant classes. $SpO_2$ is a

297     continuous variable and the use of three class is probably insufficient, especially when high

298     $SpO_2$ range is suggested as potentially harmful [18, 19]. Instead of the classification model,

299     the next step could be to test regression models' performance. $SpO_2$ was predicted at 5min

300     after ventilator setting change, a clinically relevant delay. However, the delay between

301     ventilator setting change and oxygenation steady state is not well defined and vary from 1 to

302     71 minutes according to the parameter set ($FiO_2$, PEEP or other parameters that change

303     mean airway pressure) and clinical conditions studied [15, 20, 21]. This needs further

304     research and probably more sophisticated clinical decision support systems using machine

305     learning predictive models should consider these factors. Finally, we excluded hemodynamic

306     unstable patients using a treatment criteria ($\geq$ 2 vasoactive drugs infused) because this

307     condition decreases pulse oximeter reliability [22, 23]. The validation and electronic

308     availability of reliable markers of hemodynamic instability in children such as

309     plethysmographic variability indices could be helpful [24].

310

311     # Conclusion

14

312    This pilot study using machine learning predictive model resulted in an algorithm with good

313    accuracy. To obtain a robust algorithm with such a method, more data rows are needed,

314    suggesting the need of multicenter pediatric intensive care high resolution databases.

315

# Acknowledgments

322

# References

324

325    1.    Rose L, Schultz M, Cardwell C, Jouvet P, McAuley D, Blackwood B. Automated versus non-

326    automated weaning for reducing the duration of mechanical ventilation for critically ill adults and

327    children: a cochrane systematic review and meta-analysis. Crit Care. 2015;19:48. doi:

328    10.1186/s13054-015-0755-6.

329    2.    Jouvet P, Eddington A, Payen V, Bordessoule A, Emeriaud G, Gasco R, et al. A pilot

330    prospective study on closed loop controlled ventilation and oxygenation in ventilated children during

331    the weaning phase. Crit Care. 2012;16(3):R85. doi: 10.1186/cc11343.

332    3.    Flechelles O, Ho A, Hernert P, Emeriaud G, Zaglam N, Cheriet F, et al. Simulations for

333    mechanical ventilation in children: review and future prospects. Crit Care Res Pract.

334    2013;2013:943281. doi: 10.1155/2013/943281.

335   4.      Brossier D, El Taani R, Sauthier M, Roumeliotis N, Emeriaud G, Jouvet P. Creating a High-

336   Frequency Electronic Database in the PICU: The Perpetual Patient. Pediatr Crit Care Med.

337   2018;19(4):e189-e98. doi: 10.1097/PCC.0000000000001460.

338   5.      Cakar N, Tuŏrul M, Demirarslan A, Nahum A, Adams A, Akýncý O, et al. Time required for

339   partial pressure of arterial oxygen equilibration during mechanical ventilation after a step change in

340   fractional inspired oxygen concentration. Intens Care Med 2001;27(4):655-9.

341   6.      Pediatric Acute Lung Injury Consensus Conference G. Pediatric acute respiratory distress

342   syndrome: consensus recommendations from the Pediatric Acute Lung Injury Consensus Conference.

343   Pediatr Crit Care Med. 2015;16(5):428-39. doi: 10.1097/PCC.0000000000000350.

344   7.      Les recommandations des experts de la SRLF. Le monitorage et les alarmes ventilatoires des

345   malades ventilés artificiellement. Réanim Urgences. 2000;9:407-12.

346   8.      Chawla N, Japkowicz N, A. Kotcz A. Editorial: special issue on learning from imbalanced data

347   sets. ACM SIGKDD Explorations Newsletter. 2004;6:1-6.

348   9.      Elhassan T, Aljurf M, Al-Mohanna F, Shoukri M. Classification of Imbalance Data using Tomek

349   Link (T-Link) Combined with Random Under-sampling (RUS) as a Data Reduction Method. Journal of

350   Informatics and Data Mining. 2016;1:1-12.

351   10.     Chawla N, Bowyer K, Hall L, Kegelmeyer W. SMOTE: Synthetic Minority Over-sampling

352   Technique. Journal of Artificial Intelligence Research. 2002;16:321-57. doi: 10.1613/jair.953.

353   11.     Gnana Sheela K, Deepa S. Review on methods to fix number of hidden neurons in neural

354   networks. Mathematical Problems in Engineering. 2013;2013:11. doi: 10.1155/2013/425740.425740.

355   12.     Breiman L. Bagging predictors. Berkeley: University of California, Statistics Do; 1994 421.

356   13.     Safavian S, Landgrebe D. A survey of decision tree classifier methodology. IEEE Transactions

357   on Systems, Man, and Cybernetics. 1991;21(3):660-74. doi: 10.1109/21.97458.

358   14.     Sokolova M, Lapalme G. A systematic analysis of performance measures for classification

359   tasks. Information Processing & Management. 2009;45(4):427-37. doi: 10.1016/j.ipm.2009.03.002.

360  15.    Smallwood CD, Walsh BK, Arnold JH, Gouldstone A. Equilibration Time Required for

361  Respiratory System Compliance and Oxygenation Response Following Changes in Positive End-

362  Expiratory Pressure in Mechanically Ventilated Children. Crit Care Med. 2018;46(5):e375-e9. doi:

363  10.1097/CCM.0000000000003001.

364  16.    Johnson AE, Pollard TJ, Shen L, Lehman LW, Feng M, Ghassemi M, et al. MIMIC-III, a freely

365  accessible critical care database. Sci Data. 2016;3:160035. doi: 10.1038/sdata.2016.35.

366  17.    Johnson AE, Ghassemi MM, Nemati S, Niehaus KE, Clifton DA, Clifford GD. Machine Learning

367  and Decision Support in Critical Care. Proceedings of the IEEE Institute of Electrical and Electronics

368  Engineers. 2016;104(2):444-66. doi: 10.1109/JPROC.2015.2501978.

369  18.    Girardis M, Busani S, Damiani E, Donati A, Rinaldi L, Marudi A, et al. Effect of Conservative vs

370  Conventional Oxygen Therapy on Mortality Among Patients in an Intensive Care Unit: The Oxygen-

371  ICU Randomized Clinical Trial. JAMA. 2016;316(15):1583-9. doi: 10.1001/jama.2016.11993.

372  19.    Pannu SR, Dziadzko MA, Gajic O. How Much Oxygen? Oxygen Titration Goals during

373  Mechanical Ventilation. Am J Respir Crit Care Med. 2016;193(1):4-5. doi: 10.1164/rccm.201509-

374  1810ED.

375  20.    Tugrul S, Cakar N, Akinci O, Ozcan PE, Disci R, Esen F, et al. Time required for equilibration of

376  arterial oxygen pressure after setting optimal positive end-expiratory pressure in acute respiratory

377  distress syndrome. Crit Care Med. 2005;33(5):995-1000.

378  21.    Fildissis G, Katostaras T, Moles A, Katsaros A, Myrianthefs P, Brokalaki H, et al. Oxygenation

379  equilibration time after alteration of inspired oxygen in critically ill patients. Heart Lung.

380  2010;39(2):147-52. doi: 10.1016/j.hrtlng.2009.06.009.

381  22.    Salyer J. Neonatal and pediatric pulse oximetry. Respir care. 2003;48(4):386-96.

382  23.    Fouzas S, Priftis KN, Anthracopoulos MB. Pulse oximetry in pediatric practice. Pediatrics.
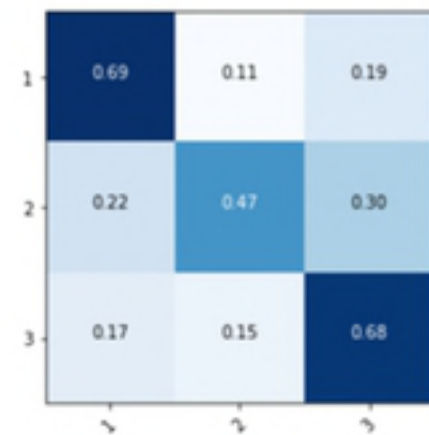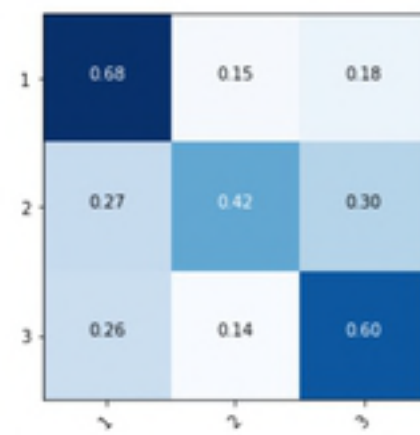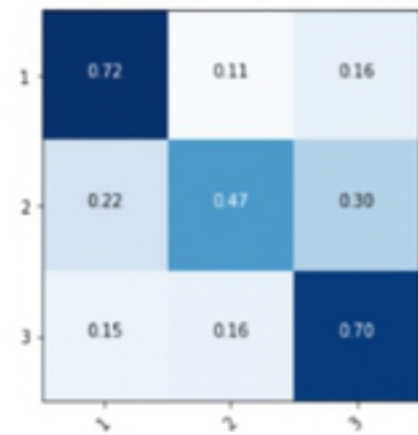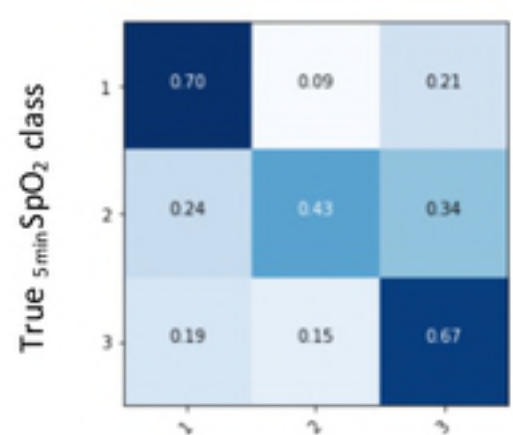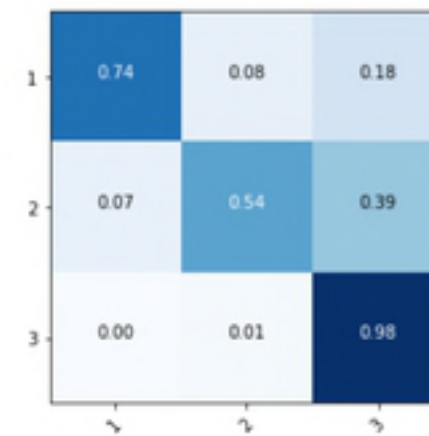
383  2011;128(4):740-52. doi: 10.1542/peds.2011-0271.
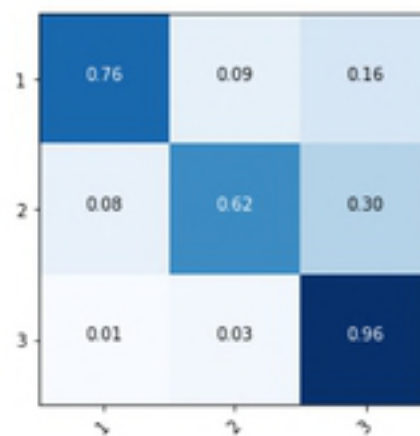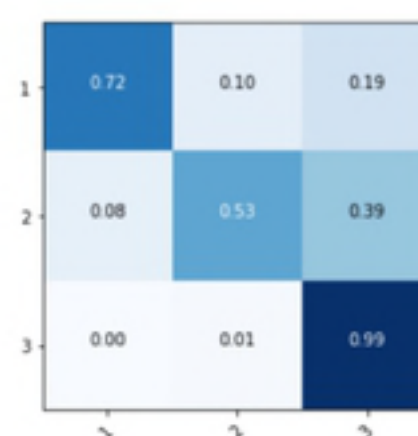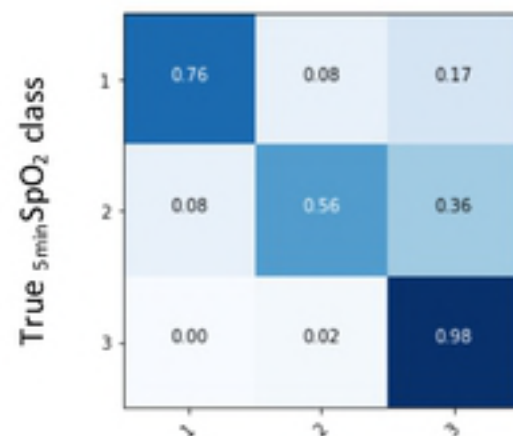
384    24.    Chandler JR, Cooke E, Petersen C, Karlen W, Froese N, Lim J, et al. Pulse oximeter

385    plethysmograph variation and its relationship to the arterial waveform in mechanically ventilated

386    children. J Clin Monit Comput. 2012;26(3):145-51. doi: 10.1007/s10877-012-9347-z.

387

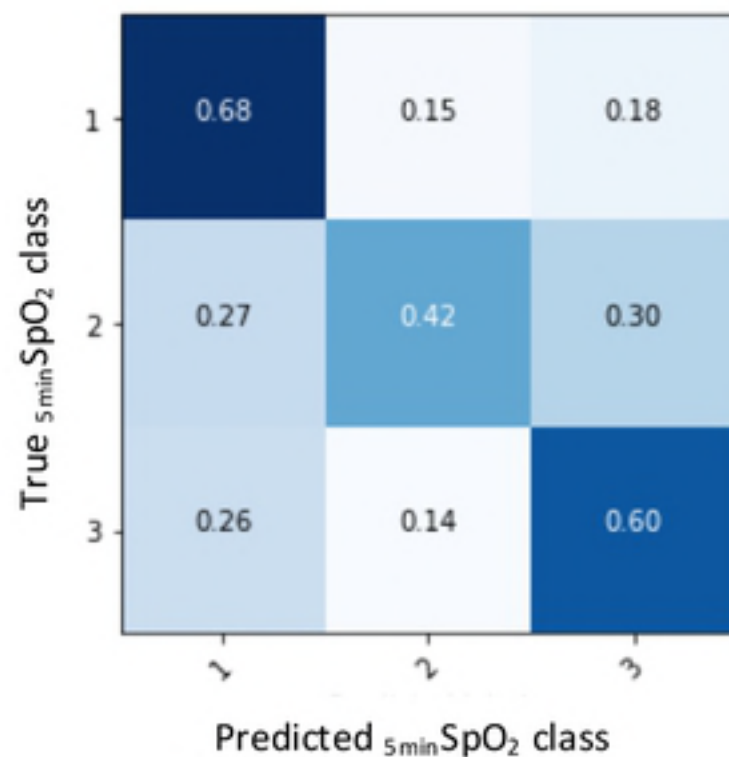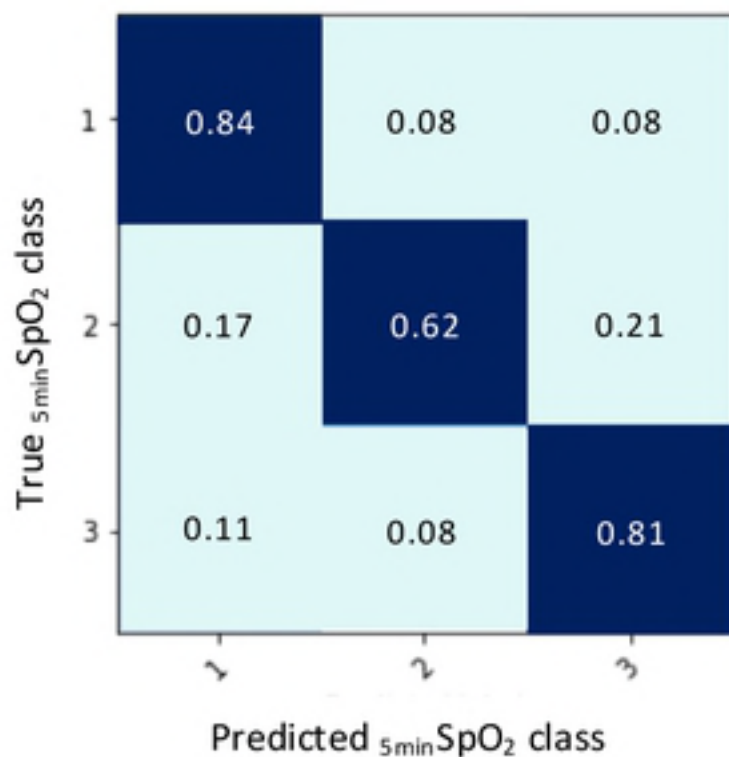388    # Supporting information

389
390    **S1 File: Data formatting process**

**ANN**



**BACDT**

Predicted $_{5min}$SpO$_2$ class

|   | A | B | C | D |

Training confusion matrix — Test confusion matrix

ANN

True $_{5min}$SpO$_2$ class / Predicted $_{5min}$SpO$_2$ class

Training (ANN):
| | 1 | 2 | 3 |
|---|---|---|---|
| 1 | 0.84 | 0.08 | 0.08 |
| 2 | 0.17 | 0.62 | 0.21 |
| 3 | 0.11 | 0.08 | 0.81 |

Test (ANN):
| | 1 | 2 | 3 |
|---|---|---|---|
| 1 | 0.68 | 0.15 | 0.18 |
| 2 | 0.27 | 0.42 | 0.30 |
| 3 | 0.26 | 0.14 | 0.60 |

BACDT

Training (BACDT):
| | 1 | 2 | 3 |
|---|---|---|---|
| 1 | >0.99 | <0.01 | <0.01 |
| 2 | 0.01 | 0.99 | 0.01 |
| 3 | <0.01 | 0.01 | 0.99 |

Test (BACDT):
| | 1 | 2 | 3 |
|---|---|---|---|
| 1 | 0.76 | 0.09 | 0.16 |
| 2 | 0.08 | 0.62 | 0.30 |
| 3 | 0.01 | 0.03 | 0.96 |

| EMR | Age |
| | Weight |

**DATABASE**

**Monitor (0.2Hz)**
- Heart
- Pulse
- SpO$_2$

**Ventilator (0.03Hz)**

Ventilator's settings
- FiO$_2$
- PEEP
- Vt
- PS above PEEP
- PC above PEEP

Ventilator's measures
- Expiratory Minute Volume
- I:E Ratio
- Measured RR
- Mean Airway Pressure
- PIP

Data formatting process

Computed, during data formatting process, from changes in FiO$_2$, PEEP & Tidal Volume
- Delta FiO$_2$
- Delta PEEP
- Delta Vt, PS or PC

ML classifiers (ANN and BACDT) training & testing process

Observed $_{5min}$SpO$_2$

Predicted $_{5min}$SpO$_2$

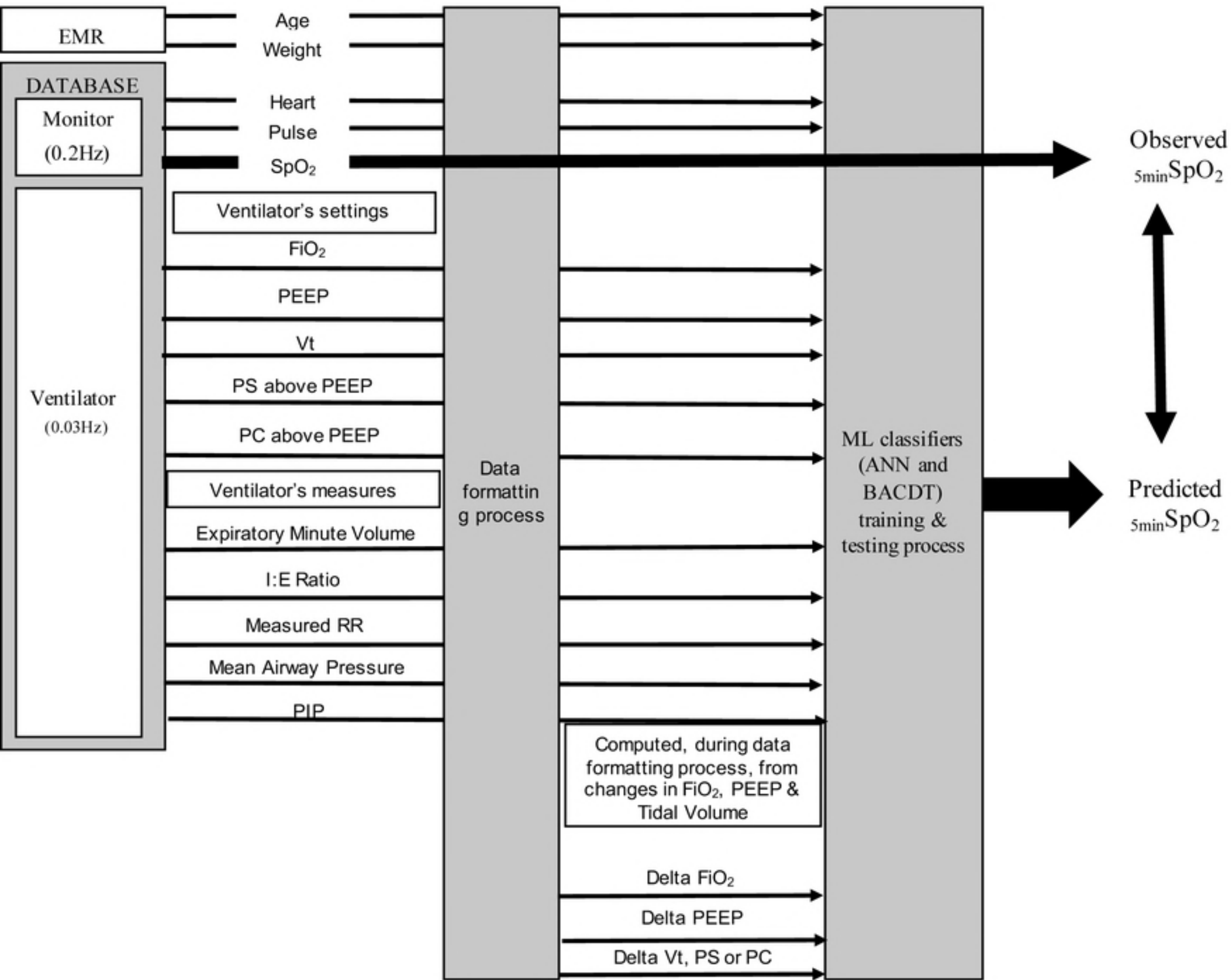| | | | |
|---|---|---|---|
| Training set: 975,036 samples<br>Test set: 193,528 samples<br>Class Balancing: TOMEK applied to dataset (before dataset has been split into training & test set) to remove tomek links, random undersampling applied to class 3 once dataset is split into training and testing sub-sets, then SMOTE applied to classes 1 and 2 to make their cardinalities equal to that of class 3 (n=325,012). | Training set: 2,293,119 samples<br>Test set: 201,926 samples<br>Class Balancing: SMOTE applied to classes 1 & 2 to make their cardinalities equal to that of class 3 (n=764,373). | Training set: 487,464 samples<br>Test set: 106,028 samples<br>Class Balancing: TOMEK applied to dataset (before dataset has been split into training & test set) to remove tomek links, random undersampling applied to class 3 once dataset is split into training and testing sub-sets, then SMOTE applied to classes 1 and 2 to make their cardinalities equal to that of class 3 (n=162,488). | Training set: 1,462,503 samples<br>Test set: 281,028 samples<br>Class Balancing TOMEK applied to dataset (before dataset has been split into training & test set) to remove tomek links, random undersampling applied to class 3 once dataset is split into training and testing sub-sets, then SMOTE applied to classes 1 and 2 to make their cardinalities equal to that of class 3 (n=487,501). |

- •• The data-set is first divided into two parts; the training-set and the test-set.
- •• The training of the "Bagged" Complex Trees includes a k-fold cross-validation, which is performed as follows:
  - ➤ Randomly partition the data-set into k equal-sized subsets (folds).
  - ➤ For each of the k equal-sized subsets:
  - ✓ Train/fit the model on the elements contained in the other (k-1) subsets.
  - ✓ Test the model's accuracy on the given subset.
  - ➤ Iterate over the k subsets, until each one has been used once for testing the model's performance during its training.
  - ➤ The training validation score consists of the average score obtained by validating the model on all k subsets.